



## King's Research Portal

DOI:

[10.1016/j.tcs.2016.05.028](https://doi.org/10.1016/j.tcs.2016.05.028)

*Document Version*

Peer reviewed version

[Link to publication record in King's Research Portal](#)

*Citation for published version (APA):*

Daykin, J. W., Groult, R., Guesnet, Y., Lecroq, T., Lefebvre, A., Léonard, M., & Prieur-Gaston, É. (2016). Binary block order Rouen transform. *Theoretical Computer Science*. <https://doi.org/10.1016/j.tcs.2016.05.028>

### **Citing this paper**

Please note that where the full-text provided on King's Research Portal is the Author Accepted Manuscript or Post-Print version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version for pagination, volume/issue, and date of publication details. And where the final published version is provided on the Research Portal, if citing you are again advised to check the publisher's website for any subsequent corrections.

### **General rights**

Copyright and moral rights for the publications made accessible in the Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognize and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Research Portal

### **Take down policy**

If you believe that this document breaches copyright please contact [librarypure@kcl.ac.uk](mailto:librarypure@kcl.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.

# Accepted Manuscript

Binary block order Rouen transform

Jacqueline W. Daykin, Richard Groult, Yannick Guesnet, Thierry Lecroq, Arnaud Lefebvre et al.

PII: S0304-3975(16)30164-5  
DOI: <http://dx.doi.org/10.1016/j.tcs.2016.05.028>  
Reference: TCS 10777

To appear in: *Theoretical Computer Science*

Received date: 7 August 2015  
Revised date: 3 May 2016  
Accepted date: 19 May 2016

Please cite this article in press as: J.W. Daykin et al., Binary block order Rouen transform, *Theoret. Comput. Sci.* (2016), <http://dx.doi.org/10.1016/j.tcs.2016.05.028>

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



## Highlights

- Novel twin binary Burrows-Wheeler type transforms are introduced.
- The transforms are defined for Lyndon-like B-words which apply binary block order.
- We call this approach the B-BWT Rouen Transform.
- These bijective Rouen transforms and inverses are computed in linear time.
- Preliminary experimental results indicate potential value of binary transforms.

# Binary Block Order Rouen Transform<sup>☆</sup>

Jacqueline W. Daykin<sup>a,b,d</sup>, Richard Groult<sup>c,d</sup>, Yannick Guesnet<sup>d</sup>, Thierry Lecroq<sup>d</sup>, Arnaud Lefebvre<sup>d</sup>, Martine Léonard<sup>d</sup>, Élise Prieur-Gaston<sup>d</sup>

<sup>a</sup>Department of Computer Science, Royal Holloway, University of London, UK

<sup>b</sup>Department of Informatics, King's College London, UK

<sup>c</sup>Modélisation, Information et Systèmes (MIS), Université de Picardie Jules Verne, Amiens, France

<sup>d</sup>Normandie Université, Université de Rouen-Normandie, LITIS EA 4108, Normastic FR CNRS 3638, 76821 Mont-Saint-Aignan, France

---

## Abstract

We introduce bijective Burrows-Wheeler type transforms for binary strings<sup>1</sup>. The original method by Burrows and Wheeler [BW94] is based on lexicographic order for general alphabets, and the transform is defined to be the last column of the ordered BWT matrix. This new approach applies binary block order,  $B$ -order, which yields not one, but twin transforms: one based on Lyndon words, the other on a repetition of Lyndon words. These binary  $B$ -BWT transforms are constructed here for  $B$ -words, analogous structures to Lyndon words. A key computation in the transforms is the application of a linear-time suffix-sorting technique, such as [KA03, KS03, KSB06, NZC09], to sort the cyclic rotations of a binary input string into their  $B$ -order. Moreover, like the original lexicographic transform, we show that computing the  $B$ -BWT inverses is also achieved in linear time by using straightforward combinatorial arguments.

*Keywords:* algorithm, bijective, binary alphabet, block order, Burrows-Wheeler Transform,  $B$ -word, data clustering, inverse transform, lexicographic order, linear, Lyndon word, string, suffix array, suffix-sorting, word

---

## 1. Introduction

The Burrows-Wheeler Transform (BWT) is a well-studied and useful text transformation scheme which was introduced by Burrows and Wheeler in 1994 [BW94]. The transform is constructed by arranging all cyclic rotations, or conjugates, of a string into a matrix in lexicographic order and extracting the rightmost column – hence a permutation of the input. Notably the transform

---

<sup>☆</sup> This paper is in honour of eminent stringologist Professor Bill Smyth: our colleague, friend and fellow life traveller. Bill is an inspiration to us all!

<sup>1</sup>These twin transforms originated in problem solving sessions of the Rouen 2012 String-Masters workshop, hence the name Rouen Transform.

has the property that it tends to cluster occurrences of the same letter in the input string into repetitions (runs), and moreover, that it is reversible. Due to this data-clustering property, the Burrows-Wheeler Transform finds wide-ranging applications in: preprocessing for data compression (it is also called block-sorting compression) and it is core to the bzip2 family of text compressors; bioinformatics and sequence processing [BCRS13, MRRS05]; and text indexing [FM00] – indeed the versatility of the transform is increasing to include, for instance, multimedia information retrieval [ABM08].

During the last decade, further research into the combinatorial properties of the BWT has led to related observations and efficient time and space implementations, including [CDP05, CGKL13, Kär07, KKP12, SLLM09], along with a range of bijective variants such as the sort transform [GS12, Kuf09].

In this paper we continue more recent research into Burrows-Wheeler type transforms based on non-lexicographic orders, thus contributing towards a taxonomy of BWT-type transforms. In [DS14] the  $V$ -BWT transform is constructed which is derived from  $V$ -order on an arbitrary alphabet; in [DW14] the  $D$ -BWT degenerate transform is introduced for indeterminate strings, namely strings of subsets of an alphabet. We propose here the binary  $B$ -BWT Rouen Transform, based on binary block  $B$ -order, to achieve text transformations for binary strings: preliminary experiments indicate potential gains of the binary Rouen Transform over the original lexicographic BWT, encouraging further practical research.

Space saving techniques with the classic BWT have been achieved by first factoring the input string into Lyndon words in linear time [CFL58, Day11, Duv83, Lot83, Smy03]. Similarly, in Algorithm 1, we efficiently factor the input string into  $B$ -words, namely analogues of Lyndon words defined using  $B$ -order [DD96, DD03] – we achieve transforms here for the case of  $B$ -words.

To sort the rotations of a binary string into a BWT matrix in  $B$ -order, we will apply a linear-time suffix-sorting algorithm, such as that of Kärkkäinen and Sanders [KS03], Kärkkäinen, Sanders and Burkhardt [KSB06], Ko and Aluru [KA03], or Nong, Zhang, and Chan [NZC09].

To our knowledge these are the first Burrows-Wheeler type transforms defined specifically for a binary alphabet; furthermore, the ordering technique applied,  $B$ -order, is based on blocks or runs of bits, and so is itself tuned to clustering – see Definition 2.1. However, it is interesting to note that the concept of modifying order was proposed in [CT98], soon after the transform was first introduced. In particular, they showed that re-ordering the characters of the alphabet provides superior compression than using their original order: the main difference with the original order is the grouping of the vowels at the beginning of the alphabet. Further, the gain for this alphabet re-ordering algorithm (with respect to the original order) for a BWT-based compressor was later shown in [AT05].

### 1.1. Notation

We introduce some standard notation; for further automata and stringology terminology, theory, algorithms and data structures see [CHL07, Smy03].

A **string**, or equivalently **word**, is a sequence of elements, letters or symbols over a set  $\Sigma$ , an **alphabet**. Throughout this paper the given alphabet is assumed to be binary,  $\Sigma = \{0, 1\}$ , and the binary elements are **bits**; the binary strings may be encoded as integers thus inducing the alphabet  $\Sigma' = \{0, 1, 2, \dots\}$ . All strings are written in mathbold:  $\mathbf{x}$ ,  $\mathbf{w}$ , and so on;  $\epsilon$  denotes the empty string.

Let  $\mathbf{x}$  be the binary string  $\mathbf{x} = \mathbf{x}[1 \dots n] = x_1 x_2 \dots x_n$ , where each  $x_i \in \{0, 1\}$ , and the **length**  $|\mathbf{x}| = n$ . Hence the  $i$ -th symbol of a string  $\mathbf{x}$  is denoted by  $\mathbf{x}[i]$ , or simply  $x_i$ . We denote by  $\mathbf{x}[i \dots j]$ , or  $x_i \dots x_j$ , the **substring** of  $\mathbf{x}$  that starts at position  $i$  and ends at position  $j$ . A string  $\mathbf{w}$  is a substring, or **factor**, of  $\mathbf{x}$  if  $\mathbf{x} = \mathbf{u}\mathbf{w}\mathbf{v}$ , where  $\mathbf{u}, \mathbf{v} \in \Sigma^*$ ; specifically, a string  $\mathbf{w} = w_1 \dots w_m$  is a substring of  $\mathbf{x} = x_1 \dots x_n$  if  $w_1 \dots w_m = x_i \dots x_{i+m-1}$  for some  $i$ . Words  $\mathbf{x}[1 \dots i]$  are called **prefixes** of  $\mathbf{x}$ , and words  $\mathbf{x}[i \dots n]$  are called **suffixes** of  $\mathbf{x}$ . The prefix  $\mathbf{u}$  (respectively suffix  $\mathbf{v}$ ) is a **proper** prefix (suffix) of a word  $\mathbf{x} = \mathbf{u}\mathbf{v}$  if  $\mathbf{x} \neq \mathbf{u}$  ( $\neq \mathbf{v}$ ).

If we can write  $\mathbf{x} = \mathbf{u}\mathbf{v} = \mathbf{w}\mathbf{u}$  for some nonempty  $\mathbf{u}$ , then we say that  $\mathbf{x}$  has **border**  $\mathbf{u}$ ; if no such  $\mathbf{u}$  exists, then  $\mathbf{x}$  is said to be **border-free** (or **unbordered**). If  $\mathbf{x} = \mathbf{u}^k$  for some nonempty  $\mathbf{u}$  and some integer  $k > 1$ , we say that  $\mathbf{x}$  is a **repetition** with **seed**  $\mathbf{u}$ ; otherwise, we say that  $\mathbf{x}$  is **primitive**. If a string  $\mathbf{x} = \mathbf{u}\mathbf{v}$ , then  $\mathbf{v}\mathbf{u}$  is said to be a **rotation** (cyclic shift) or **conjugate** of  $\mathbf{x}$ .

The subject of this paper requires careful study of rotations of strings with associated notation. For a binary string  $\mathbf{x}$ , we let  $\mathbf{x}^+$ ,  $\mathbf{x}_0^+$ ,  $\mathbf{x}_1^+$  each denote a clockwise rotation of (any) one bit, a 0 bit, and a 1 bit respectively. Likewise, let  $\mathbf{x}^-$ ,  $\mathbf{x}_0^-$ ,  $\mathbf{x}_1^-$  denote an anti-clockwise rotation of (any) one bit, a 0 bit, and a 1 bit respectively. For sequences of bit rotations, let  $(\mathbf{x}^+)^k$  denote  $k$  clockwise rotations, and similarly for  $(\mathbf{x}_0^-)^k$  and so on. For instance, if  $\mathbf{x} = 1010010$ , then  $\mathbf{x}^+ = 0100101$  and  $\mathbf{x}^- = 0101001$ ; further, since it is known that the rightmost bit of  $\mathbf{x}$  is 0, then we can apply  $\mathbf{x}_0^-$ , giving  $\mathbf{x}_0^- = \mathbf{x}^- = 0101001$ . Unless specified otherwise, a rotation will be assumed to be anti-clockwise.

We now identify four disjoint sets of rotations of  $\mathbf{x}$  which arise in the  $B$ -BWT binary transform algorithms:

- $\mathcal{R}_0^u$  - unbordered rotations commencing with 0 and ending with 1
- $\mathcal{R}_1^u$  - unbordered rotations commencing with 1 and ending with 0
- $\mathcal{R}_0^b$  - bordered rotations commencing with 0 and ending with 0
- $\mathcal{R}_1^b$  - bordered rotations commencing with 1 and ending with 1.

Note that with our focus on rotating single bits, we use a weak form of border-free here: for instance, if only the first and last bits are known, the string  $001 \dots 001$  would be considered border-free.

We may encode a non-empty binary string  $\mathbf{x}$  as the string of integers  $\hat{\mathbf{x}} = Z_0 N_1 Z_1 N_2 Z_2 \dots N_r Z_r$ , where  $Z_0$  is the number of 0's at the start of  $\mathbf{x}$ , then  $N_1$  is the number of following 1's, and so on. For example, if  $\mathbf{x} = 0111001111$

then  $r = 2$  and  $\hat{x} = 13240$ . In general  $Z_0, Z_r \geq 0$  but  $N_1, Z_1, \dots, N_r \geq 1$ . Given an encoding  $\hat{x} = Z_0 N_1 Z_1 N_2 Z_2 \dots N_r Z_r$ , let  $\hat{N}(\mathbf{x}) = \hat{N} = N_1 N_2 \dots N_r$ , and similarly  $\hat{Z}(\mathbf{x}) = \hat{Z} = Z_0 Z_1 Z_2 \dots Z_r$ . Notation such as  $N_i(\mathbf{x})$  (respectively  $Z_i(\mathbf{x})$ ) would indicate the  $i$ th  $N$  (respectively  $Z$ ) value of an encoded string  $\hat{x}$ . The maximal substrings  $0^{Z_i}$  and  $1^{N_j}$  of  $\mathbf{x}$  are called **blocks**, or in stringology they are usually called **runs**; the terminology  $N$  block ( $Z$  block) will mean a block of 1's (0's).

Given the encoding  $\hat{x} = Z_0 N_1 Z_1 N_2 Z_2 \dots N_r Z_r$  of a border-free string  $\mathbf{x}$ , if a rotation(s) of  $\mathbf{x}$  causes  $\mathbf{x}$  to become bordered, that is by rotating the bits of a prefix/suffix of a block, then that block becomes a **split block** (otherwise a **non-split** or **complete block**); for instance, if  $\mathbf{x} = 010010111$ , where  $\hat{x} = Z_0 N_1 Z_1 N_2 Z_2 N_3 Z_3 = 1121130$ , then  $(x_1^-)^2 = 110100101$  and the block associated with  $N_3 = 3$  has become a split block. We call maximal substrings in  $\mathbf{x}$  of the form  $1^{N_i} 0^{Z_i}$  a **block pair**.

Let the **weight**  $\omega \mathbf{x}$  of a binary string  $\mathbf{x}$  be the total number of 1's in  $\mathbf{x}$ , that is,  $\omega \mathbf{x} = x_1 + x_2 + \dots + x_n$ ; thus all rotations in the set of cyclic rotations, or conjugates, of  $\mathbf{x}$  have the same weight  $\omega$ .

This paper is also concerned with methods for totally ordering strings. If a string  $\mathbf{x}$  is less than a string  $\mathbf{y}$  in **lexorder** (lexicographic order), we write  $\mathbf{x} < \mathbf{y}$  and  $\mathbf{y} > \mathbf{x}$ . The next section discusses a binary string ordering technique.

## 2. Binary Block Order ( $B$ -order)

We start with the basic definitions and computational issues for working with  $B$ -order in Section 2.1, followed by some combinatorial results in Section 2.2 which will be applied when computing a transform and inverse.

### 2.1. Block Order Preliminaries

We introduce here the method for totally ordering binary strings which will be used in Sections 3 & 4 for constructing the  $B$ -BWT matrices followed by extracting the transforms.

**Definition 2.1.** [DD96]  $B$ -order (**Block total order**  $\ll$ ) of words of length  $n$ . Given  $\mathbf{x} = x_1 x_2 \dots x_n$  with encoding  $\hat{x} = Z_0 N_1 Z_1 N_2 Z_2 \dots N_r Z_r$  and  $\mathbf{y} = y_1 y_2 \dots y_n$  with encoding  $\hat{y} = Z'_0 N'_1 Z'_1 N'_2 Z'_2 \dots N'_s Z'_s$ , where each  $x_i, y_j \in \{0, 1\}$  and  $|\mathbf{y}| = |\mathbf{x}|$ , then  $\mathbf{x} \ll \mathbf{y}$  if

- (i)  $\omega \mathbf{x} < \omega \mathbf{y}$ , or
- (ii)  $\omega \mathbf{x} = \omega \mathbf{y}$  and there is a least  $i$  in  $1 \leq i \leq \min\{r, s\}$  with  $N_i \neq N'_i$  and then  $N_i > N'_i$ , or
- (iii)  $\omega \mathbf{x} = \omega \mathbf{y}$  and  $r = s$  and  $N_i = N'_i$  for  $1 \leq i \leq r$ , and there is a least  $j$  in  $0 \leq j < r$  with  $Z_j \neq Z'_j$  and then  $Z_j < Z'_j$ .

That is, in  $B$ -order, strings are ordered first by weight, then by lexorder on the sizes of blocks of 1's (over the alphabet  $\{1 > 2 > 3 > \dots\}$ ) and finally by lexorder on the sizes of blocks of 0's (over the alphabet  $\{0 < 1 < 2 < 3 < \dots\}$ ).

We will describe sorting using this ordering technique as *B-sorting* or *B-ordering* strings.

A primitive operation is the comparison of two strings  $\mathbf{x}, \mathbf{y}$  in *B-order*. This can be achieved by straightforward bit comparisons; similarly, following the linear computation of  $\hat{\mathbf{x}}, \hat{\mathbf{y}}$ , for instance by Run Length Encoding, applying linear lexordering to the weights  $\omega$ , the  $\hat{N}$  values and the  $\hat{Z}$  values as necessary, we get:

**Claim 2.2.** *The comparison of two strings  $\mathbf{x}, \mathbf{y}$  in *B-order*, that is to determine  $\mathbf{x} \ll \mathbf{y}$ , can be achieved in linear time.*

In this paper we will demonstrate twin transforms for the case of *B-words* which are defined as follows:

**Definition 2.3.** [DD03] *A binary primitive string  $\mathbf{x}$  is a *B-word* if it is the unique minimum in *B-order* among the cyclic rotations of  $\mathbf{x}$ .*

Hence *B-words* are analogues of Lyndon words – specified as lexicographically least in their set of conjugates [CFL58, Lot83] – while defined with *B-order* rather than the original lexorder; another analogue of Lyndon words, defined with *V-order*, is given in [DD03, DS14].

**Example 2.4.** *Consider the cyclic rotations of the string 110100. In *B-order*:  $110100 \ll 011010 \ll 001101 \ll 100110 \ll 010011 \ll 101001$ . Therefore the given string 110100 is a *B-word*.*

Furthermore: 0, 1, and strings (block pairs) of the form  $1^j 0^k$  with  $j, k > 0$  are *B-words*; any other string  $\mathbf{x}$  which is a *B-word* is in the set  $\mathcal{R}_1^u$  with  $Z_0(\mathbf{x}) = 0$ . Note that *B-words* share similar properties to Lyndon words such as being primitive, border-free and must start with the largest  $N$  value; hence we denote  $N_1$  as *max*. We now give a characterization:

So let  $\mathbf{x} = x_1 x_2 \dots x_n \in \mathcal{R}_1^u$  be a *B-word*; then  $\hat{\mathbf{x}} = Z_0 N_1 Z_1 N_2 Z_2 \dots N_r Z_r$  with  $Z_0(\mathbf{x}) = 0$ . If  $\hat{N}(\mathbf{x}) = (N_1 \dots N_p)^k$  then let  $\zeta(\mathbf{x}) = \zeta_1 \zeta_2 \dots \zeta_k$  with  $\zeta_i = Z_{(i-1)p+1} Z_{(i-1)p+2} \dots Z_{(i-1)p+p}$ . Notation such as  $\zeta_i(\mathbf{x})$  would indicate the  $i$ th  $\zeta$  value of an encoded string  $\hat{\mathbf{x}}$ .

**Theorem 2.5.** *A string  $\mathbf{x} \in \mathcal{R}_1^u$  is a *B-word* if and only if either  $\hat{N}(\mathbf{x})$  forms a Lyndon word over the alphabet  $\{1 > 2 > 3 > \dots\}$ , or,  $\hat{N}(\mathbf{x})$  forms a repetition of a Lyndon word over the alphabet  $\{1 > 2 > 3 > \dots\}$  and  $\zeta(\mathbf{x})$  forms a Lyndon word over the alphabet  $\{1^p < 1^{p-1} 2 < 1^{p-1} 3 < \dots\}$ .*

*Proof.* ( $\Rightarrow$ ) Suppose first that  $\mathbf{x}$  is a *B-word*. Then Definition 2.1 ensures that  $\hat{N}(\mathbf{x})$  is either a Lyndon word over  $\{1 > 2 > 3 > \dots\}$  (i) or a repetition of Lyndon words over the alphabet  $\{1 > 2 > 3 > \dots\}$ .

Consider the case that  $\hat{N}(\mathbf{x})$  is a repetition of Lyndon words over  $\{1 > 2 > 3 > \dots\}$ . If  $\zeta(\mathbf{x})$  does not form a Lyndon word over the alphabet  $\{1^p < 1^{p-1} 2 < 1^{p-1} 3 < \dots\}$  then there exists a rotation  $\mathbf{r}$  of  $\mathbf{x}$  such that  $\zeta(\mathbf{r}) \ll \zeta(\mathbf{x})$



which contradicts Definition 2.3 stating that  $\mathbf{x}$  is strictly smaller than all of its rotations.

( $\Leftarrow$ ) Suppose now that  $\hat{N}(\mathbf{x})$  and  $\zeta(\mathbf{x})$  satisfy either of the stated Lyndon properties for a string  $\mathbf{x} \in \mathcal{R}_1^u$ . Then by Definition 2.1,  $\mathbf{x}$  must precede any of its rotations in  $B$ -order, including those with split blocks, and hence forms a  $B$ -word.  $\square$

**Example 2.6.** Let  $\mathbf{x} = 11101001110010$  then  $\hat{N}(\mathbf{x}) = 3131 = (31)^2$ ,  $\hat{Z}(\mathbf{x}) = 01221$ ,  $\zeta(\mathbf{x}) = 12 \cdot 21$  with  $\zeta_1(\mathbf{x}) = 12 < 21 = \zeta_2(\mathbf{x})$ :  $\hat{N}(\mathbf{x})$  is a repetition of Lyndon words and  $\zeta(\mathbf{x})$  is a Lyndon word over their respective alphabets thus  $\mathbf{x}$  is a  $B$ -word. Now let  $\mathbf{x}' = 11100101110100$  which is a conjugate of  $\mathbf{x}$  then  $\hat{N}(\mathbf{x}') = 3131 = (31)^2$ ,  $\hat{Z}(\mathbf{x}') = 02112$ ,  $\zeta(\mathbf{x}') = 21 \cdot 12$  with  $\zeta_1(\mathbf{x}') = 21 > 12 = \zeta_2(\mathbf{x}')$ :  $\hat{N}(\mathbf{x}')$  is a repetition of Lyndon words but  $\zeta(\mathbf{x}')$  is not a Lyndon word over their respective alphabets thus  $\mathbf{x}'$  is not a  $B$ -word. Note that the border-free property of the Lyndon word  $\zeta(\mathbf{x})$  is given by the words  $\zeta_i$  and not their individual letters.

We will refer to the first case in Theorem 2.5 as **type primitive  $B$ -words**,  **$B$ -prim**, and the second as **type repetition  $B$ -words**,  **$B$ -rep** – these two types of words yield two different methods for computing and inverting transforms to recover the input.

We will require all  $B$ -ordered rotations of a string for computing a transform, illustrated by:

**Example 2.7.** Consider the  $B$ -BWT matrix  $M_{\mathbf{x}}$  of cyclic rotations for  $\mathbf{x} = 110110010$ , ordered into block order, giving  $M_{\mathbf{x}}^B$ :

1	1	0	1	1	0	0	1	0
0	1	1	0	1	1	0	0	1
1	1	0	0	1	0	1	1	0
0	1	1	0	0	1	0	1	1
1	0	1	1	0	1	1	0	0
0	1	0	1	1	0	1	1	0
0	0	1	0	1	1	0	1	1
1	0	1	1	0	0	1	0	1
1	0	0	1	0	1	1	0	1

In this example with  $\mathbf{x}$  a  $B$ -prim, we have  $\hat{N}(\mathbf{x}) = 221$ , a Lyndon word (using  $>$  rather than the usual  $<$ ); hence, by Definition 2.1(ii), all rotations with non-split  $N$  blocks occur after  $\mathbf{x}$  in  $M_{\mathbf{x}}$ . Again by Definition 2.1(ii), rotations in  $M_{\mathbf{x}}$  with split  $N$  blocks will also occur after  $\mathbf{x}$ . Note that  $\hat{Z}(\mathbf{x}) = 0121$  is not then required to be a Lyndon word (although it is using  $<$ ), since the  $B$ -word criteria of  $\mathbf{x}$  is guaranteed here by the Lyndon property of  $\hat{N}$ .

Note that to test if a string  $\mathbf{x}$  which starts with 1 is a  $B$ -word, by Lyndon properties, we need only look at rotations of  $\mathbf{x}$  in the set  $\mathcal{R}_1^u$ , as any rotation

starting with a  $Z \neq 0$  block has a prefix of 0's of length greater than the required  $Z_0(\mathbf{x}) = 0$ , satisfying Definition 2.1(iii). Similarly, for rotations of  $\mathbf{x}$  in the set  $\mathcal{R}_1^b$  which arise from split  $N$  blocks, Definition 2.1(ii) applies.

Furthermore, we can factor an arbitrary binary string  $\mathbf{x}$  of length  $n$  bits into  $B$ -words in linear time using a modification of Duval's  $O(n)$  Lyndon decomposition algorithm [Duv83, Day11]: first compute  $\hat{\mathbf{x}}$  in linear time using Run Length Encoding; then, in tandem, factor the  $N$ 's into Lyndon words using  $>$ , while for any detected repetitions of Lyndon words, factor the  $Z$  components into Lyndon words using  $<$ , taking the  $B$ -words to be those factors with maximal length – see Algorithm 1. So, without incurring any computational inefficiency, we will assume that the input to the transform is a  $B$ -word.

---

**Algorithm 1** Factor a binary string  $\mathbf{x}$  into  $B$ -words in  $O(n)$ -time.

---

```

procedure COMPUTE-BF( $\mathbf{x}$ )
   $F \leftarrow \emptyset$ ;
  Compute  $\hat{\mathbf{x}} = Z_0 N_1 Z_1 N_2 Z_2 \dots N_r Z_r$ ;
  ▷ compute the  $\hat{\mathbf{x}}$  encoding using Run Length Encoding

  Compute  $NLF(N_1 N_2 \dots N_r) = \mathbf{n}_1 \mathbf{n}_2 \dots \mathbf{n}_s$ ;
  ▷ compute the Lyndon factorization of  $N$  components using  $>$ 

  for each repetition  $(\mathbf{n}_p = \mathbf{n}_{p+1} = \dots = \mathbf{n}_q)$  detected during factoring do
    Compute  $\zeta LF(\mathbf{n}_p = \mathbf{n}_{p+1} = \dots = \mathbf{n}_q) = \theta_1 \theta_2 \dots \theta_t$ ;
    ▷ compute Lyndon factorization of  $Z$  components corresponding
    ▷ to  $(\mathbf{n}_p = \dots = \mathbf{n}_q)$  using  $<$ 

    for  $i \leftarrow 1$  to  $t$  do
      add end index of  $\theta_i$  to  $F$ 
    end for
  end for
  return  $F$ 
end procedure

```

---

Note that due to the maximality property of the Lyndon factorization, the factorization into  $B$ -words will likewise be maximal and thus unique. To formalize this, we apply the following lemma with the terminology: an *UMFF* is a unique maximal factorization family *FF*, where an *FF* is a subset  $\mathcal{W} \subseteq \Sigma^+$  which permits, for every nonempty string  $\mathbf{x}$  on  $\Sigma$ , a factorization of  $\mathbf{x}$  over  $\mathcal{W}$ .

**Lemma 2.8.** (*The  $\mathbf{xyz}$  Lemma [DD03]*). *An FF  $\mathcal{W}$  is an UMFF if and only if whenever  $\mathbf{xy}, \mathbf{yz} \in \mathcal{W}$  for some nonempty  $\mathbf{y}$ , then  $\mathbf{xyz} \in \mathcal{W}$ .*

**Lemma 2.9.** *The set of  $B$ -words forms an UMFF.*

*Proof.* Let  $\mathbf{x}, \mathbf{y}, \mathbf{z}$  be three non-empty words on  $\Sigma$  such that  $\mathbf{xy}$  and  $\mathbf{yz}$  are  $B$ -words. We will use the fact that Lyndon words form an UMFF [CFL58,

DD03, Day11] and therefore satisfy Lemma 2.8. Since  $B$ -words must start with a 1, we have  $\mathbf{xy}[1] = \mathbf{yz}[1] = 1$  and therefore the overlap  $\mathbf{y}[1] = 1$ . We proceed to consider four cases arising from Theorem 2.5, and we assume  $\mathbf{x}$  and  $\mathbf{z}$  are nonempty for otherwise it is trivial. We will also use the result that a word is a Lyndon word if and only if it is lexicographically less than any of its nonempty proper suffixes [Duv83, Smy03]; additionally, if  $\mathbf{u}, \mathbf{v}$  are Lyndon words then  $\mathbf{uv}$  or  $\mathbf{vu}$  is a Lyndon word.

- i.  $\hat{N}(\mathbf{xy})$  and  $\hat{N}(\mathbf{yz})$  form Lyndon words over  $\{1 > 2 > 3 > \dots\}$ . Since Lyndon words form an UMFF,  $\hat{N}(\mathbf{xyz})$  is a Lyndon word over  $\{1 > 2 > 3 > \dots\}$ .
- ii.  $\hat{N}(\mathbf{xy})$  and  $\hat{N}(\mathbf{yz})$  form repetitions of Lyndon words over  $\{1 > 2 > 3 > \dots\}$ , and  $\hat{Z}(\mathbf{xy})$  and  $\hat{Z}(\mathbf{yz})$  form Lyndon words over  $\{0 < 1 < 2 < 3 < \dots\}$ . Suppose that  $\hat{N}(\mathbf{xy}) = \mathbf{n}^p$  and  $\hat{N}(\mathbf{xyz}) = \mathbf{n}^q$ , where  $p < q$ , then  $\zeta(\mathbf{xyz})$  is a Lyndon word over the alphabet  $\{1^{|\mathbf{n}|} < 1^{|\mathbf{n}|-1}2 < 1^{|\mathbf{n}|-1}3 < \dots\}$ . Otherwise,  $\hat{N}(\mathbf{xyz})$  is not a repetition, and consider the seed  $\mathbf{n}$  of  $\hat{N}(\mathbf{xy})$  and the seed  $\mathbf{m}$  of  $\hat{N}(\mathbf{yz}) = \mathbf{m}^r$ . Let  $i > 1$  be the index such that  $\hat{N}(\mathbf{y})[1] = \hat{N}(\mathbf{xy})[i]$ ; then there are two cases for the overlap  $\hat{N}(\mathbf{y})$ .

If  $i$  indexes the start of some  $\mathbf{n}$  in  $\mathbf{n}^p$  then, by the maximality of Lyndon words,  $\mathbf{m}$  cannot be shorter than  $\mathbf{n}$ . If  $\mathbf{m}$  is the same length as  $\mathbf{n}$ , then  $\hat{N}(\mathbf{xyz})$  forms a repetition as above; otherwise the suffix of  $\hat{N}(\mathbf{xy})$  starting at  $i$  is a proper prefix of  $\mathbf{m}$  and  $\hat{N}(\mathbf{xyz})$  is a Lyndon word over  $\{1 > 2 > 3 > \dots\}$ .

If on the other hand a prefix  $\hat{N}(\mathbf{v})$  of  $\hat{N}(\mathbf{y})$  is a proper suffix of some  $\mathbf{n}$ , then, since  $\hat{N}(\mathbf{n}) < \hat{N}(\mathbf{v})$ , it must be the  $p$ th  $\mathbf{n}$ . Writing  $\mathbf{n} = \mathbf{uv}$ , and applying Lemma 2.8 with nonempty overlap  $\hat{N}(\mathbf{v})$ , we have  $\hat{N}(\mathbf{um})$  is a Lyndon word over  $\{1 > 2 > 3 > \dots\}$  with prefix  $\mathbf{n}$ , and likewise  $\hat{N}(\mathbf{n}^{p-1}\mathbf{um})$  and  $\hat{N}(\mathbf{n}^{p-1}\mathbf{um}^r)$ . Therefore  $\hat{N}(\mathbf{xyz})$  forms a Lyndon word as required.

- iii.  $\hat{N}(\mathbf{xy})$  satisfies Case [i] and  $\hat{N}(\mathbf{yz})$  satisfies Case [ii]. Let  $\hat{N}(\mathbf{y}) = \mathbf{m}_1\mathbf{m}_2\dots\mathbf{m}_s[1\dots t]$ , where  $\mathbf{m}_1 = \mathbf{m}_2 = \dots = \mathbf{m}_s = \mathbf{m}$  and  $s < r$ ,  $t \leq |\mathbf{m}|$  or  $s \leq r$ ,  $t < |\mathbf{m}|$ . If  $t = |\mathbf{m}|$ , then since  $\mathbf{m}$  is a suffix of  $\hat{N}(\mathbf{xy})$  we have  $\hat{N}(\mathbf{xy}) < \hat{N}(\mathbf{m})$ , and it follows that  $\hat{N}(\mathbf{xym}^{r-s})$  is a Lyndon word. If  $t < |\mathbf{m}|$ , then write  $\mathbf{m} = \mathbf{uv}$  where  $\mathbf{u} = \mathbf{m}_s[1\dots t]$  and  $\mathbf{v} = \mathbf{m}_s[t+1\dots|\mathbf{m}|]$ . Applying Lemma 2.8 with nonempty overlap  $\hat{N}(\mathbf{u})$ , we have that  $\hat{N}(\mathbf{xym})$  is a Lyndon word with suffix  $\hat{N}(\mathbf{m})$ . Furthermore,  $\hat{N}(\mathbf{xym}^{r-s})$  is a Lyndon word. In either case  $\hat{N}(\mathbf{xyz})$  forms a Lyndon word over  $\{1 > 2 > 3 > \dots\}$ .
- iv.  $\hat{N}(\mathbf{xy})$  satisfies Case [ii] and  $\hat{N}(\mathbf{yz})$  satisfies Case [i]. Consider the prefix of  $\hat{N}(\mathbf{yz})$  which overlaps the  $i$ th repetition  $\mathbf{n}_i = \mathbf{n}$  in  $\hat{N}(\mathbf{xy})$  with minimal  $i$ , that is  $\hat{N}(\mathbf{yz})[1\dots\hat{N}(\mathbf{xy})[i|\mathbf{n}|]] = \hat{N}(\mathbf{xy})[d\dots e] = \mathbf{v}$ . So we can write  $\mathbf{n}_i = \mathbf{uv}$ ,  $\mathbf{v} \neq \varepsilon$ .

If  $d$  indexes the start of the repetition  $\mathbf{n}_i$ , that is  $d = (i-1)|\mathbf{n}| + 1$ , then the Lyndon word  $\hat{N}(\mathbf{yz})$  has prefix  $\hat{N}(\mathbf{n}_i)$  and so  $\hat{N}(\mathbf{n}^{i-1}\mathbf{yz})$  is also a Lyndon word.

Otherwise,  $d$  does not index the start of  $\mathbf{n}_i$ . Since  $\hat{N}(\mathbf{n}_i)$  and  $\hat{N}(\mathbf{yz})$  are both Lyndon words with nonempty overlap  $\hat{N}(\mathbf{v})$ , then Lemma 2.8 applies, and  $\hat{N}(\mathbf{xyz})$  is a Lyndon word with prefix  $\hat{N}(\mathbf{n}_i)$ . As above  $\hat{N}(\mathbf{n}^{i-1}\mathbf{xyz})$  is also a Lyndon word.

In either case  $\hat{N}(\mathbf{xyz})$  forms a Lyndon word over  $\{1 > 2 > 3 > \dots\}$ .

We conclude that the set of  $B$ -words satisfies the xyz Lemma and hence forms an UMFF.  $\square$

For further results, examples and discussion on  $B$ -words, see [DD96, DD03].

## 2.2. Combinatorics of Rotations in $B$ -order

Analogously to the classic lexorder BWT [BW94], the method for a binary block order transform,  $B$ -BWT, is specified as: order the conjugates of an input string  $\mathbf{x}$  of length  $n$  into  $B$ -order, giving the  $n \times n$   $B$ -BWT matrix  $M$ ; the **transform**  $T$  is then defined to be the last right-hand column of the matrix  $M$ , that is  $M[i, n]$  for  $1 \leq i \leq n$  – in Example 2.7 the  $B$ -BWT transform is  $T = 010100111$ . The particular cases considered here are where  $\mathbf{x}$  is assumed to be a  $B$ -word of either type; hence we work with Definition 2.1(ii) & (iii) (part (i) does not apply to cyclic rotations).

From Definition 2.1(ii), we see that the sizes of the left-most  $N_1$  components of rows in  $M$  are decreasing. This partitions  $M$  into **groups**  $G_t, G_{t-1}, \dots, G_1$  of rows where group  $G_i$  is all rows where the first occurrence of a block of 1's is of size  $i$ ,  $N_1(\mathbf{x}) \geq i \geq 1$ . Hence in Example 2.7,  $M$  has groups:  $G_2$ , rows 1 - 4; and  $G_1$ , rows 5 - 9.

We proceed to establish some new results for  $B$ -words, first for primitive then general and finally type repetition, which although fairly simple, help to unravel patterns in the maze of bits which lead to the transforms and inverses.

The next lemma shows that rotations of the  $Z$  suffixes of  $B$ -prim words which end with 0 can just be inserted directly into the matrix  $M$ .

**Lemma 2.10.** *Let the binary string  $\mathbf{r}$  be a row in the  $B$ -BWT matrix  $M$  of a  $B$ -prim word with  $\hat{\mathbf{r}} = Z_0 N_1 Z_1 N_2 Z_2 \dots N_r Z_r$  and  $0 < Z_r = k$ . Then  $\mathbf{r} \ll (\mathbf{r}^-)^1 \ll \dots \ll (\mathbf{r}^-)^k$  in  $M$ ; furthermore, these rotations are consecutive in  $M$ .*

*Proof.* Definition 2.1(ii) shows that rows in  $M$  with the same  $N$  values will be consecutive in  $M$ . Rotating a 0 suffix bit leaves the  $N$  values of a string unchanged but increments the  $Z_0$  component – hence Definition 2.1(iii) applies.  $\square$

**Corollary 2.11.** *Let  $\mathbf{r}$  be as in Lemma 2.10. Then the sequence of rotations  $\mathbf{r}, (\mathbf{r}^-)^1, \dots, (\mathbf{r}^-)^k$  corresponds to a substring of the form  $0^k$  in the transform  $T$ .*

That is, each  $Z$  block in a given  $B$ -prim word  $\mathbf{x}$  will correspond to a  $Z$  suffix of some row in  $M$ , which by the subsequent rotations, will give a run of  $0^k$  in the transform  $T$  – a key insight for computing the inverse of the transform.

The  $N$  blocks of unit size, which form group  $G_1$ , also have particular importance for initiating the inverse transforms; the next corollary addresses the formation of this group.

**Corollary 2.12.** *Following the rotation of a complete  $Z$  block suffix in (either type)  $M$ , rotating the next adjacent 1 bit gives a row in group  $G_1$  of  $M$ .*

Hence rotating the right-most bit of any  $N$  block (of any size) in the input string  $\mathbf{x}$  generates a rotation in group  $G_1$ . We also observe that in any row of the matrix  $M$ , at most one block of either type  $N$  or  $Z$  can be split, thus giving a bordered row in  $\mathcal{R}_1^b$  or  $\mathcal{R}_0^b$  respectively.

The following lemma will be applied later to show that, given a partial  $B$ -BWT matrix of a  $B$ -prim word derived from a suffix array, the full  $B$ -BWT matrix can be readily constructed.

**Lemma 2.13.** *Let  $\mathbf{r}_i, \mathbf{r}_j$  be two rows in a  $B$ -BWT matrix  $M$  with  $\mathbf{r}_i \ll \mathbf{r}_j$ . Suppose that  $\mathbf{r}_i[n] = \mathbf{r}_j[n]$ , and also that these rows do not satisfy  $\mathbf{r}_i[1] = 0$  and  $\mathbf{r}_j[1] = 1$ . Then  $\mathbf{r}_i^- \ll \mathbf{r}_j^-$ .*

*Proof.* Since the strings  $\mathbf{r}_i, \mathbf{r}_j$  are conjugates, they cannot differ only on their right-hand blocks of bits, and so the ordering of  $\mathbf{r}_i$  and  $\mathbf{r}_j$  is determined prior to the last blocks. Consider the cases for the first bits:

- $\mathbf{r}_i[1] = \mathbf{r}_j[1]$ . Then clearly Definition 2.1(ii) & (iii) apply to each of the eight cases (start 0 or 1; end 0 or 1; differ on 0's or 1's) as appropriate, giving  $\mathbf{r}_i^- \ll \mathbf{r}_j^-$ .
- $\mathbf{r}_i[1] = 1$  and  $\mathbf{r}_j[1] = 0$ . If  $\mathbf{r}_i[n] = \mathbf{r}_j[n] = 1$ , then  $N_1((\mathbf{r}_i)_1^-) > N_1((\mathbf{r}_j)_1^-) = 1$ , and so by Definition 2.1(ii),  $\mathbf{r}_i^- \ll \mathbf{r}_j^-$ . Otherwise  $\mathbf{r}_i[n] = \mathbf{r}_j[n] = 0$ , whereby  $1 = Z_0((\mathbf{r}_i)_0^-) < Z_0((\mathbf{r}_j)_0^-)$ , and so either the order of  $\mathbf{r}_i^-$  and  $\mathbf{r}_j^-$  is determined by the  $N$  components which haven't changed, or the order is determined by Definition 2.1(iii), and in either case  $\mathbf{r}_i^- \ll \mathbf{r}_j^-$ .

□

In other words, if two strings in the matrix both end with either 0 or 1, but do not start with 0 and 1 respectively, then the  $B$ -order of the strings is maintained when rotating those last bits – note that this does not apply to the case of different end bits. In particular, this lemma leads the way to applying existing knowledge on the order of strings to the ordering of their rotations.

The Last First mapping for the classic Burrows Wheeler transform utilizes the fact that the relative positions of letters correspond in the first and last columns in the BWT matrix – hence it is applied for recovering a string from the computed transform. Lemma 2.13 shows that this correspondence is not

the case with binary bits: see Example 2.7 where the 1st, 2nd, 3rd, 4th and 5th 1 bits in the Last column map to the 3rd, 4th, 5th, 1st and 2nd 1 bits in the First column respectively. Hence we distinguish between matching suffix block pairs of the form  $10^i$  with  $i > 0$  and with  $i = 0$  in the mapping in Section 3.4.

**Lemma 2.14.** *Suppose that the binary strings  $\mathbf{u} = 1 \dots 10^i$  and  $\mathbf{v} = 1 \dots 10^j$  where  $\mathbf{u} \ll \mathbf{v}$  are conjugates of one  $B$ -prim word. If  $i, j > 0$  or  $i, j = 0$ , then  $(\mathbf{u}^-)^{i+1} \ll (\mathbf{v}^-)^{j+1}$ .*

*Proof.* Consider the two cases.

**Case  $i, j > 0$ .** Let  $\hat{\mathbf{u}} = N_1 Z_1 N_2 Z_2 \dots N_r Z_r$  and  $\hat{\mathbf{v}} = N'_1 Z'_1 N'_2 Z'_2 \dots N'_s Z'_s$ . Since  $\mathbf{u}, \mathbf{v}$  are unbordered conjugates of a Lyndon word, then Definition 2.1(iii) cannot apply; furthermore,  $r = s$ . Then  $\hat{N}(\mathbf{u}) \neq \hat{N}(\mathbf{v})$  and by Definition 2.1(ii) some  $N_t > N'_t$ . Further, since  $\mathbf{u}, \mathbf{v}$  are conjugates with the same weight, they cannot differ only on the last  $N_r, N'_s$  pair. So assume  $t < r$ , and then the order of  $\mathbf{u}$  or  $\mathbf{v}$  is not changed by the rotations of the respective suffixes  $10^i$  and  $10^j$ . In this case, rotations of suffixes always map to group  $G_1$  - taking you down the matrix.

**Case  $i, j = 0$ .** Here  $\mathbf{u}, \mathbf{v} \in \mathcal{R}_1^b$ , and let  $\hat{\mathbf{u}} = N_1 Z_1 N_2 Z_2 \dots N_{r+1}$  and  $\hat{\mathbf{v}} = N'_1 Z'_1 N'_2 Z'_2 \dots N'_{s+1}$ ; then  $r+1 = s+1$ . Suppose Definition 2.1(iii) applies, then  $N_1 + N_{r+1} = N'_1 + N'_{s+1}$  and  $N_q = N'_q$  for  $2 \leq q \leq r$ , contradicting the primitive property of the  $B$ -prim word. Therefore Definition 2.1(ii) applies as above. In this case, the rotation of a suffix bit always increments the group size of each conjugate - taking you up the matrix, possibly to different groups.

□

Since Lemma 2.14 holds for all such  $\mathbf{u}, \mathbf{v}$ , then we see that the  $k$ th occurrence of a suffix  $10^i$  in Last maps to the  $k$ th occurrence of  $10^i$  as a prefix in group  $G_1$ , and the  $k$ th occurrence of a row belonging to  $\mathcal{R}_1^b$  in  $G_h$  maps to the  $k$ th occurrence of the prefix  $1^{h+1}$  in group  $G_{h+1}$ .

We conclude these observations on rotations with a result for the repetition case of Theorem 2.5, which shows that rotations of  $B$ -rep words with the same  $\hat{N}$  encodings, arising from repetitions, are adjacent in the  $B$ -BWT matrix  $M$ .

**Lemma 2.15.** *Suppose that the rotations  $\mathbf{r}_i, \mathbf{r}_{i+1}, \dots, \mathbf{r}_{i+j}$  of a  $B$ -rep word  $\mathbf{x}$  belong to the set  $\mathcal{R}_1^u$  and satisfy  $\hat{N}(\mathbf{r}_i) = \hat{N}(\mathbf{r}_{i+1}) = \dots = \hat{N}(\mathbf{r}_{i+j})$ . Then these  $j+1$  rows are consecutive in the  $B$ -BWT matrix  $M$ ; furthermore, the rotations of their  $Z$  suffixes all immediately follow  $\mathbf{r}_i, \mathbf{r}_{i+1}, \dots, \mathbf{r}_{i+j}$  in the matrix  $M$ .*

*Proof.* Suppose that the conjugate  $\mathbf{q}$  of  $\mathbf{x}$  has identical  $N$  components to  $\mathbf{r}_i$ , that is  $\hat{N}(\mathbf{q}) = \hat{N}(\mathbf{r}_i)$ , and also that  $\mathbf{q}$  has a non-empty  $Z$  prefix, so  $Z_0(\mathbf{q}) > 0$ . Since the rotations  $\mathbf{r}_i, \dots, \mathbf{r}_{i+j}$ , which are generated by repetitions, all have empty  $Z$  prefixes ( $Z_0 = 0$ ), then, by Definition 2.1(iii), these  $j+1$  rows must

all precede  $\mathbf{q}$  in  $M$  – the next conjugate adjacent to  $\mathbf{r}_{i+j}$  in  $M$  must have the same  $\hat{N}$  encoding and  $Z_0 = 1$ , followed by those with  $Z_0 \geq 1$  prefixes but all with the same  $\hat{N}$  encoding. Suppose now that  $\mathbf{q}$  has distinct  $N$  components to  $\mathbf{r}_i$ , that is  $\hat{N}(\mathbf{q}) \neq \hat{N}(\mathbf{r}_i)$ . Then by Definition 2.1(ii),  $\mathbf{q}$  either precedes  $\mathbf{r}_i$ , or comes after  $\mathbf{r}_i, \dots, \mathbf{r}_{i+j}$  and all the rotations of their  $Z$  suffixes.  $\square$

We are now ready to introduce the binary block order Burrows-Wheeler type transforms. For  $B$ -prim words, the transform is presented in Section 3, followed by the computation of its inverse in Section 3.1; for  $B$ -rep words, the transform is given in Section 4 and the inverse in Section 4.1.

### 3. The $B$ -BWT Binary Transform - Primitive Case

The key computation in the  $B$ -BWT transform is the  $B$ -ordering of all cyclic rotations of an input  $\mathbf{x}$  to obtain the  $B$ -BWT matrix  $M$  – the transform  $T$  is then readily extracted from the right-hand column of  $M$ . For the  $B$ -sorting we apply a linear time and space suffix array construction for sorting string suffixes and hence conjugates – we will basically apply such a method as a “black box”; for details the interested reader is referred to Kärkkäinen and Sanders [KS03], Ko and Aluru [KA03], or Nong, Zhang and Chan [NZC09].

These suffix array techniques apply lexorder to compare string suffixes which can be readily adapted to binary  $B$ -order comparison – Claim 2.2 shows that the linear complexity is maintained. Note that for the modification for the  $B$ -order suffix array, the suffixes relate to conjugates which all have the same weight  $\omega$ , therefore part (i) of Definition 2.1 can be ignored and parts (ii) & (iii) adapted and applied as follows. Suppose that  $\mathbf{u}$  and  $\mathbf{v}$  are nonempty suffixes such that  $|\mathbf{u}| < |\mathbf{v}|$ . If some  $\hat{N}_i(\mathbf{u}) \neq \hat{N}_i(\mathbf{v})$ , then the order is determined by the two  $N_i$  values according to Definition 2.1(ii). If  $\hat{N}(\mathbf{u}) = \hat{N}(\mathbf{v}) \neq \varepsilon$ , then  $\hat{Z}_0(\mathbf{u}) \neq \hat{Z}_0(\mathbf{v})$ , and the order is determined according to Definition 2.1(iii). If  $\hat{N}(\mathbf{u})$  is empty, then  $\mathbf{u}$  is a single  $Z$  block, and so we append  $1^{N_1(\mathbf{x})}$  as a suffix to  $\mathbf{u}$ ; similarly for  $\hat{N}(\mathbf{v})$ . For cases when  $\hat{N}(\mathbf{u})$  is a proper prefix of  $\hat{N}(\mathbf{v})$  we apply:

**Lemma 3.1.** *Let  $\mathbf{b}$  and  $\mathbf{bd}$  be distinct nonempty proper suffixes of a  $B$ -prim word  $\mathbf{x}$ , where  $\mathbf{b}$  starts with a suffix of the  $j^{\text{th}}$   $N$  block corresponding to  $\hat{N}_j(\mathbf{x})$ , and  $\mathbf{bd}$  starts with a suffix of the  $i^{\text{th}}$   $N$  block corresponding to  $\hat{N}_i(\mathbf{x})$ ,  $i < j$ . With  $\mathbf{x} = \mathbf{ab} = \mathbf{cbd}$ , where  $\mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{d} \neq \varepsilon$ , then the conjugates  $\mathbf{ba}$  and  $\mathbf{bdc}$  of  $\mathbf{x}$  satisfy  $\mathbf{ba} \ll \mathbf{bdc}$ .*

*Proof.* Note that since  $\hat{N}(\mathbf{x})$  is a Lyndon word, and hence border-free, then  $\hat{N}(\mathbf{b})$  cannot also be a prefix of  $\hat{N}(\mathbf{x})$ . Further,  $|\hat{N}_1(\mathbf{bd})| \leq |\hat{N}_i(\mathbf{x})|$  and  $|\hat{N}_1(\mathbf{b})| \leq |\hat{N}_j(\mathbf{x})|$ .

Consider  $B$ -ordering the conjugates  $\mathbf{ba}$  and  $\mathbf{bdc}$  by first applying lexorder (with  $>$ ) to  $\hat{N}(\mathbf{ba})$  and  $\hat{N}(\mathbf{bdc})$ . Since the prefixes of the conjugates are both  $\mathbf{b}$ , and  $\mathbf{c} \neq \varepsilon$ , then  $|\mathbf{a}| = |\mathbf{dc}|$  and  $|\mathbf{a}| > |\mathbf{d}|$ . From the Lyndon properties of  $\hat{N}(\mathbf{x})$  we have  $\hat{N}(\mathbf{x}) \ll \hat{N}(\mathbf{d})$ , but since  $\hat{N}(\mathbf{x})$  is border-free  $\hat{N}(\mathbf{d})$  is not a prefix of  $\hat{N}(\mathbf{a})$ . Hence, by minimality of  $\mathbf{x}$ ,  $\hat{N}(\mathbf{a})[h] > \hat{N}(\mathbf{d})[h]$  for some  $1 \leq h \leq |\mathbf{d}|$ ,

and therefore the order is not changed by appending  $c$  to  $d$ ; we conclude that  $ba \ll bdc$ .  $\square$

Hence, although suffixes relating to conjugates may not have the same length, we can order them similarly to lexorder by using their integer  $\hat{N}$  and  $\hat{Z}$  values – intermediary values relating to split blocks can be recorded during a right-to-left Run Length Encoding scan. Depending on the implementation, the  $B$ -order suffix array can be constructed just for the  $N$  values, and then the  $Z$  prefixes and suffixes entered into  $M$  according to Lemma 2.10. Overall, by modifying the ordering for a classic suffix array technique we get:

**Claim 3.2.** *The  $B$ -BWT matrix  $M$  can be computed in time linear in the bit length,  $n$ , of the input  $x$ .*

Finally, the transform  $T$  is trivially extracted from  $M$  as  $M[i, n]$ ,  $1 \leq i \leq n$ .

**Example 3.3.** *Consider the  $B$ -prim word  $x = 1110011001110101100$  where  $\hat{x} = 03222311122$ . Then the  $n \times n$  matrix of  $B$ -sorted rotations is given by  $M_x^B$ :*

1	1	1	0	0	1	1	0	0	1	1	1	0	1	0	1	1	0	0
0	1	1	1	0	0	1	1	0	0	1	1	1	0	1	0	1	1	0
0	0	1	1	1	0	0	1	1	0	0	1	1	1	0	1	0	1	1
1	1	1	0	1	0	1	1	0	0	1	1	1	0	0	1	1	0	0
0	1	1	1	0	1	0	1	1	0	0	1	1	1	0	0	1	1	0
0	0	1	1	1	0	1	0	1	1	0	0	1	1	1	0	0	1	1
1	1	0	0	1	1	1	0	0	1	1	0	0	1	1	1	0	1	0
0	1	1	0	0	1	1	1	0	0	1	1	0	0	1	1	1	0	1
1	1	0	0	1	1	1	0	1	0	1	1	0	0	1	1	1	0	0
0	1	1	0	0	1	1	1	0	1	0	1	1	0	0	1	1	1	0
0	0	1	1	0	0	1	1	1	0	1	0	1	1	0	0	1	1	1
1	1	0	0	1	1	0	0	1	1	1	0	1	0	1	1	0	0	1
1	1	0	0	1	0	1	1	0	0	1	1	1	0	0	1	1	0	0
1	0	0	1	1	1	0	0	1	1	0	0	1	1	1	0	1	0	1
1	0	0	1	1	1	0	1	0	1	1	0	0	1	1	1	0	0	1
1	0	1	1	0	0	1	1	1	0	0	1	1	0	0	1	1	1	0
0	1	0	1	1	0	0	1	1	1	0	0	1	1	0	0	1	1	1
1	0	0	1	1	0	0	1	1	1	0	1	0	1	1	0	0	1	1
1	0	1	0	1	1	0	0	1	1	1	0	0	1	1	0	0	1	1

The  $B$ -BWT transform is  $T = 0^2 10^2 1010^2 1^5 01^3$ .

The following small example shows better clustering of bits with the  $B$ -BWT compared to the classic transform with binary lexorder.

**Example 3.4.** *Consider the  $B$ -prim word  $x = 110100010$ . Then the  $B$ -BWT is  $010^3 101^2$  while the binary lexorder transform, with  $0 < 1$ , is  $10101010^2$  (hence  $0^2 1010101$  with  $1 < 0$ ).*



### 3.1. The Inverse B-BWT - Primitive Case

In this section we establish that the B-BWT is bijective for B-prim words by showing how to recover the B-word input losslessly from the transform.

We will assume that the given encoded B-prim word  $\hat{\mathbf{x}} = N_1 Z_1 N_2 Z_2 \dots N_r Z_r$  has  $r \geq 2$ , for otherwise the inverse is trivial:

**Lemma 3.5.** *The encoding of a B-prim word  $\mathbf{x}$  is  $\hat{\mathbf{x}} = N_1 Z_1$  if and only if the last column  $T$  of the associated B-BWT matrix  $M$  has the form  $0^{Z_1} 1^{N_1}$ .*

*Proof.* Given a B-prim word  $\mathbf{x}$  such that  $\hat{\mathbf{x}} = N_1 Z_1$ , then applying Definition 2.1: the first  $Z_1$  rows in  $M$  end with 0 due to rotations  $\mathbf{r}$  of  $\mathbf{x}$  of the form  $\mathbf{r}_0^-$ , and the ordering of  $Z$ s using  $<$ ; the following  $N_1$  rows end with 1 due to rotations of the form  $\mathbf{r}_1^-$ , and the ordering of  $N$ s using  $>$ .

Given a transform  $T = 0^{Z_1} 1^{N_1}$ , then by Corollary 2.11 there is only one  $Z$  block in  $\mathbf{x}$  of length  $Z_1$ . Since the B-word starts 1 and ends 0, then it must be  $1^{N_1} 0^{Z_1}$ , and hence type B-prim.  $\square$

### 3.2. Outline of the Inverse.

An outline of the steps for the B-BWT inversion algorithm is as follows:

- Assume that the input  $\mathbf{x}$  to be recovered is a B-prim word, and hence in the set  $\mathcal{R}_1^u$ .
- The input for the inversion is the binary transform  $T$ : the last column of bits of the B-BWT matrix  $M$ .
- Determine the number of block pairs in  $\mathbf{x}$  from  $T$  – Algorithm 2.
- Find the start index of each group in  $M$  – Algorithm 3.
- Obtain the indexes for the changes in the values of the second  $N$  blocks,  $N_2$ , in  $M$  (since  $\hat{N}(\mathbf{x})$  is primitive the  $N_2$  values cannot all be the same) – Algorithm 4.
- Apply the  $N_2$  indexes to perform a Last First matching by using the given transform  $T$  and block pairs to recover the input string – Algorithm 5.

### 3.3. Indexing for the Inverse

We now present the details of the inversion algorithm, separated into simple procedures – hence we assume the variables are global as appropriate.

In order to index the matrix  $M$ , the transform is first scanned to find the number of block pairs in the input string  $\mathbf{x}$ : applying Corollary 2.11, the number of runs of 0's in the transform gives the value of  $r$  in  $\hat{\mathbf{x}}$ , that is the number of block pairs in  $\mathbf{x}$ .

Using the linear-time procedure COMPUTE-BP, we next find the start index of each group in  $M$ : due to the splitting of  $N$  blocks there must be a group for

---

**Algorithm 2** Compute the number of block pairs in the input  $\mathbf{x}$  from  $T$ .

---

```

procedure COMPUTE-BP( $T, n$ )
     $\triangleright$  scan  $T$  and calculate the number  $r$  of block pairs in  $\mathbf{x}$ 

     $i := 1$ ;  $\triangleright$  initialize the transform index
     $blockpairs := 0$ ;  $\triangleright$  initialize the number of block pairs
    while  $i \leq n$  do
        while  $T[i] = 0$  do
             $i := i + 1$ ;
        end while
         $blockpairs := blockpairs + 1$ ;  $\triangleright$  apply Corollary 2.11
        while  $i \leq n$  and  $T[i] = 1$  do
             $i := i + 1$ ;
        end while
    end while

    return  $blockpairs$ 
end procedure

```

---

each size  $1, \dots, N_1$ . From  $r = blockpairs$ , we obtain the number,  $t$ , of rows in group  $G_1$  as follows. For each unbordered row ending  $N_i > 1$ , by rotating the suffix bit of the block  $1^{N_i}$ , the block becomes a split block, and the rotation will have the form  $10 \dots 1^{N_i-1}$  and hence is in  $\mathcal{R}_1^b$  and  $G_1$ . For any  $N_i = 1$ , the block will occur in  $\mathbf{x}$  as a substring of the form  $\dots 0^{j_i} 10^{k_i} \dots$ , then by Corollary 2.11, it is associated with the substring  $0^{k_i} 1$  in the transform  $T$ ; the last row related to this substring in  $T$  has the form  $0^{k_i} \dots 0^{j_i} 1$ . Hence  $t = r + j_i$  for all such  $i$ . Note that we can also fill in the prefixes of 0's for these rows in  $M$ .

Similarly, a rotation of each of the  $s$  bordered rows in  $G_1$  generates a row in  $G_2$ , and so we also have the number of blocks  $N_i$  where  $N_i = 2$ ; accordingly, there are  $s + l_i$  rows in group  $G_2$ , where  $l_i$  is the number of rotations associated with each  $N_i$  block occurring in  $\mathbf{x}$  as  $\dots 0^{l_i} 110^{m_i} \dots$  – and so the bootstrapping continues until the first left-hand column of  $N$ s is completed (including  $Z$  prefixes where  $Z > 0$ ). To summarize: after deducing the number of rows in group  $G_1$  from the number of block pairs, the size of the next larger  $N$  value group is obtained iteratively, moderated by the  $Z$  blocks, since their rotations do not form rows in the neighbouring group – see Algorithm 3.

Observing that the algorithm performs a single pass through the transform, with associated constant work for each bit in the transform, we get:

**Claim 3.6.** *Procedure COMPUTE-GI( $T, n$ ) runs in time linear in the bit length of the transform.*

In order to perform a Last First type mapping in Algorithm 5, we will need information on the second column of 1's in  $M$ , namely the index of each change

---

**Algorithm 3** Compute the index of each group in  $M$  from  $T$ .

---

```

procedure COMPUTE-GI( $T, n$ )
     $G[1 \dots]$ ;            $\triangleright$  declare an array list for group indexes in  $M$ 
     $j := 1$ ;              $\triangleright$  initialize the group index
     $G[j] := n$ ;           $\triangleright$  initialize the index of the start of group  $G_1$ 
     $nonsplit := 0$ ;       $\triangleright$  initialize number of non split blocks in each group
     $exp := 0$ ;           $\triangleright$  initialize the exponent for a  $Z$  block
     $bool := false$ ;      $\triangleright$  set flag for non split blocks
                         $\triangleright$   $blockpairs$  is computed in Algorithm 2

    while  $G[j] > 1$  do    $\triangleright$  find the indexes for group  $G_1$  up to group  $G_{N_1}$ 
        while  $G[j] > n - blockpairs - exp$  do
            while  $G[j] > n - blockpairs - exp$  and  $T[G[j]] = 1$  do  $\triangleright N$  suffix
                 $G[j] := G[j] - 1$ 
            end while
            while  $G[j] > n - blockpairs - exp - 1$  and  $T[G[j]] = 0$  do  $\triangleright Z$  suffix
                 $G[j] := G[j] - 1$ ;
                 $exp := exp + 1$ ;
                 $bool := true$ 
            end while
            if  $bool$  then  $nonsplit := nonsplit + 1$ ;  $bool := false$ 
            end if
        end while

         $n := G[j]$ ;        $\triangleright$  iterate through the groups
         $G[j] := G[j] + 1$ ;  $j := j + 1$ ;
         $G[j] := n$ ;        $\triangleright$  initialize the index of the start of next group
         $blockpairs := blockpairs - nonsplit$ ;  $\triangleright$  adjust for  $G_1$  (Corollary 2.12)
         $nonsplit := 0$ ;
         $exp := 0$ 
    end while

     $max := j - 1$ ;       $\triangleright max = N_1$  is used in Algorithms 4 and 5
    return  $G[1 \dots]$ 
end procedure

```

---

in  $N_2$  value. The  $B$ -BWT matrix consists of groups  $G_i$ , with  $i > 1$ , and group  $G_1$  – we use this partition to deduce the  $N_2$  column.

Procedure COMPUTE-INDEX $N_2$  first scans down the transform from  $G_{max}$  to  $G_2$  detecting  $Z$  suffixes – rotations of these suffixes generate rows in  $G_1$  by Corollary 2.12; the relevant point is that for a  $Z$  block suffix of a row in group  $G_i$  of the form  $1^i \dots 10^j$ , then  $j + 1$  rotations yield a row in  $G_1$  with prefix  $10^j 1^i 0 \dots$ , and so we know the  $N_2$  value  $i$  as it relates to a non-split block. This process gives all entries in the array  $IndexN_2$  for  $G_1$ . Using the partially completed array  $IndexN_2$ , the algorithm then sweeps back up the transform incrementally deducing the entries for larger groups. Note that although we compute a  $max \times max$  array, the number of entries is less than  $n$ , hence an  $O(n)$  method.

### 3.4. The Last First Mapping

We finally emulate the Last First mapping of the classic BWT in which the  $i$ th occurrence of a letter in the last column is mapped to the  $i$ th occurrence of that letter in the first column; however, this property only holds for 0 bits in the  $B$ -BWT matrix  $M$  (Lemma 2.10, Corollary 2.12, Lemma 2.13). Therefore, using Lemma 2.14, we apply the property to a partition of  $M$ : unbordered and bordered rows, both beginning with 1, which map to group  $G_1$  and groups  $G_i$ ,  $i > 1$ , respectively. The  $B$ -BWT Last First mapping thus rotates any suffixes of 0's until a 1 bit is reached, and then the prefix matching process is resumed, where a prefix always starts with a 1 bit.

From the first run  $0^j$  in the transform  $T$  and Corollary 2.11, we know that the input  $B$ -word  $x$  ends  $10^j$ . From Algorithm 3 we know the value  $N_1 = max$  and so, since we assume  $r \geq 2$ , the string formed by  $(x^-)^{j+1}$  has the prefix  $10^j 1^{N_1} 0$ , which, by Definition 2.1, must be the first row, rotation  $r_1$  say, in group  $G_1$ . This forms the start of the Last First pattern matching process, which proceeds as follows. Consider  $T[i]$  corresponding to  $r_1$ : suppose  $T[i] = 1$ , then by Definition 2.1,  $(r_1)^-$  is the first row in group  $G_2$  as it begins with the prefix  $110^j 1^{N_1} 0$ . Otherwise,  $T[i] = 0$ , and we perform  $(r_i)^-$  repeatedly, that is scan  $T$ , until we find a 1 bit, and hence by Corollary 2.12 the final rotation is in group  $G_1$ . The method oscillates through the groups until  $n$  bits have been recovered, thus forming a cycle.

In order to implement the Last First mapping, we need the index of the next row to obtain the prefix of the string being incrementally developed. If the next row is in group  $G_i$ , then it will come after all rows in any groups with larger  $i$  value. Furthermore, the next row will also come after any rows in  $G_i$  whose prefixes have already been selected and concatenated to the evolving string. For this we track the indexes of both  $N_1$ , with the variable  $currG$ , and  $N_2$ , with  $prevG$ , for rows as we trace the cycle of rotations to reconstruct the input  $x$ .

In the implementation in Algorithm 5 we only visit each bit in the transform once, with constant work for each step, hence:

**Claim 3.7.** *Procedure COMPUTE-LF( $T, n$ ) can be implemented in linear time.*

---

**Algorithm 4** Compute indexes for the second column of 1's,  $N_2$ , in  $M$  from  $T$ .

---

```

procedure COMPUTE-INDEXN2( $T, n$ )
   $IndexN_2[max, max]$ ;  $\triangleright$  declare an array for  $N_2$  indexes of each group
   $\triangleright$  find  $N_2$  indexes for group  $G_1$ 
   $runs := 0$ ;  $\triangleright runs$  is the number of runs of 0's in  $T$  for a group
   $IndexN_2[1, max] := G[1]$ ;  $\triangleright max$  and  $G$  are computed in Algorithm 3
  for  $i = max - 1$  downto 1 do  $\triangleright i$  indexes groups
     $j = G[i + 1]$ ;  $bool := false$ ;  $\triangleright j$  indexes rows in the current group
    repeat
      if  $T[j] = 0$  then  $bool := true$ 
      else if  $bool$  then  $runs := runs + 1$ ;  $bool := false$ 
      end if
       $j := j + 1$ 
    until  $j = G[i]$ 
     $h := IndexN_2[1, i + 1]$ ;  $N := 0$ ;  $\triangleright h$  indexes rows in group  $G_1$ 
    repeat
      if  $T[h] = 1$  then  $N := N + 1$ 
      end if  $\triangleright N$  counts  $N$  suffixes
       $h := h + 1$ 
    until  $N = runs$ 
     $IndexN_2[1, i] := h$ ;  $runs := 0$ 
  end for  $\triangleright$  iterate from  $G_1$  to  $G_{max}$  with each group bootstrapping the next
   $IndexN_2[1, 0] := n + 1$ ;  $\triangleright$  set a dummy entry for the loop
   $last := G[1] - 1$ ;  $\triangleright$  index of last row of current group
  for  $i = 1$  to  $max - 1$  do
    for  $j = 1$  to  $max$  do
      if  $last > 0$  then
         $N := 0$ ;  $bool := false$ ;  $\triangleright k$  indexes rows in previous group
        for  $k = IndexN_2[i, j]$  to  $IndexN_2[i, j - 1] - 1$  do
          if  $T[k] = 0$  then  $bool := true$ 
          else if  $bool$  then  $bool := false$ 
          else  $N := N + 1$ 
          end if
        end for
         $N_g := 0$ ;  $\triangleright N_g$  counts  $N$  suffixes in current group
        while  $N_g < N$  do
          if  $T[last] = 1$  then  $N_g := N_g + 1$ ;  $last := last - 1$ 
          end if
          while  $T[last] = 0$  do  $last := last - 1$ 
          end while
        end while
         $IndexN_2[i + 1, j] := last + 1$ 
      end if
    end for
     $IndexN_2[i + 1, 0] := IndexN_2[i, max]$ ;
  end for
  return  $IndexN_2[max, max]$ 
end procedure

```

---

---

**Algorithm 5** Compute the Last First Mapping for a transform  $T$ .

---

**procedure** COMPUTE-LF( $T, n$ )

$p := 10^{Z_r};$   $\triangleright$  initialize the pattern for the matching process

$\triangleright$  from scanning the first run of 0's in  $T$

$exp = 0;$   $\triangleright$  initialize the exponent for a  $Z$  block

$currG := 1;$   $\triangleright$  size of current group

$prevG := max;$   $\triangleright$  size of previous group from Algorithm 3

**while**  $|p| < n$  **do**  $\triangleright$  iterate until  $n$  characters have been recovered

$t := IndexN_2[currG, prevG];$   $\triangleright$  obtain index of current row in  $M$  from  
 $\triangleright IndexN_2$  computed in Algorithm 4

**if**  $T[t] = 1$  **then**  $\triangleright$  the last bit of the current rotation in  $M$  is 1

$p := 1 \circ p;$   $\triangleright$  concatenate the suffix bit 1 to the start of  $p$

$IndexN_2[currG, prevG] := IndexN_2[currG, prevG] + 1;$

$\triangleright$  increment the  $N_2$  group row index

$currG := currG + 1$   $\triangleright$  increment the group size

**else**  $\triangleright$  the last bit of the current rotation in  $M$  is 0

**repeat**  $exp := exp + 1; t := t + 1$

**until**  $T[t] = 1$

$p := 10^{exp} \circ p;$   $\triangleright$  concatenate a  $Z$  block to the start of  $p$

$IndexN_2[currG, prevG] := IndexN_2[currG, prevG] + exp + 1;$

$\triangleright$  increment the  $N_2$  group row index

$prevG := currG;$

$currG := 1;$   $\triangleright$  the prefix  $10^{exp}$  indicates a rotation in group  $G_1$  (Corollary 2.12)

$exp := 0$   $\triangleright$  reset the exponent for a  $Z$  block

**end if**

**end while**

**return**  $LF(x)$   $\triangleright$  the recovered input  $x$  is given by the pattern  $p$

**end procedure**

---

At each step of the inversion, a block pair  $10^i$ ,  $i \geq 0$ , is concatenated to the evolving string and matched to the current least unselected prefix in  $M$ , hence:

**Claim 3.8.** *Procedure COMPUTE-LF( $T, n$ ) correctly recovers the input.*

#### 4. The $B$ -BWT Binary Transform - Repetition Case

With reference to Theorem 2.5, we start with an example of a  $B$ -rep word  $\mathbf{x}$ : the encoding  $\hat{N}(\mathbf{x})$  forms a repetition of Lyndon words over  $\{1 > 2 > \dots\}$  while  $\zeta(\mathbf{x})$  is a Lyndon word over the alphabet  $\{1^p < 1^{p-1}2 < 1^{p-1}3 < \dots\}$  (where  $p$  is the length of the repeated Lyndon word in  $\hat{N}(\mathbf{x})$ ).

**Example 4.1.** *Consider the  $B$ -rep word  $\mathbf{x} = 11101100111011000$  where  $\hat{\mathbf{x}} = 031223123$ ,  $\hat{N}(\mathbf{x}) = (32)^2$  with exponent of repetition 2, and  $\hat{Z}(\mathbf{x}) = 01213$ . Then the  $n \times n$  matrix of  $B$ -sorted rotations is given by  $M_{\mathbf{x}}^B$ :*

1	1	1	0	1	1	0	0	1	1	1	0	1	1	0	0	0
1	1	1	0	1	1	0	0	0	1	1	1	0	1	1	0	0
0	1	1	1	0	1	1	0	0	1	1	1	0	1	1	0	0
0	1	1	1	0	1	1	0	0	0	1	1	1	0	1	1	0
0	0	1	1	1	0	1	1	0	0	1	1	1	0	1	1	0
0	0	1	1	1	0	1	1	0	0	0	1	1	1	0	1	1
0	0	0	1	1	1	0	1	1	0	0	1	1	1	0	1	1
1	1	0	0	1	1	1	0	1	1	0	0	0	1	1	1	0
1	1	0	0	0	1	1	1	0	1	1	0	0	1	1	1	0
0	1	1	0	0	1	1	1	0	1	1	0	0	0	1	1	1
0	1	1	0	0	0	1	1	1	0	1	1	0	0	1	1	1
1	1	0	1	1	0	0	1	1	1	0	1	1	0	0	0	1
1	1	0	1	1	0	0	0	1	1	1	0	1	1	0	0	1
1	0	0	1	1	1	0	1	1	0	0	0	1	1	1	0	1
1	0	0	0	1	1	1	0	1	1	0	0	1	1	1	0	1
1	0	1	1	0	0	1	1	1	0	1	1	0	0	0	1	1
1	0	1	1	0	0	0	1	1	1	0	1	1	0	0	1	1

The  $B$ -BWT( $\mathbf{x}$ ) =  $0^5 1^2 0^2 1^8$  with better clustering than the binary lexicographic BWT( $\mathbf{x}$ ) =  $0^2 1^2 0^2 1^4 0^2 1^2 0^1 1^2$  (with  $1 < 0$ ).

#### Example 4.2.

[i] Consider the  $B$ -rep word  $\mathbf{x} = 11101101110011000$ . Then the  $B$ -BWT( $\mathbf{x}$ ) =  $0^3 1^4 0^1 1^1 0^2 1^1 0^1 1^7$  with one higher exponent than the binary lexicographic BWT( $\mathbf{x}$ ) =  $0^3 1^2 0^1 1^4 0^1 1^2 0^2 1^2$  (with  $1 < 0$ ).

[ii] Consider the  $B$ -rep word  $\mathbf{x} = 11101011100100$ . Then the  $B$ -BWT( $\mathbf{x}$ ) =  $0^3 1^4 0^2 1^1 0^1 1^3$  with better clustering than the binary lexicographic BWT( $\mathbf{x}$ ) =  $0^2 1^2 0^1 1^1 0^1 1^1 0^1 1^2 0^1 2^2$  (with  $1 < 0$ ).

In general, for a  $B$ -rep word with repetition exponent  $k$  and  $s$  distinct  $N$  components, the  $B$ -BWT matrix  $M$  partitions as follows ( $Z_{l,r}$  denotes the  $Z_r$

suffix of row  $l$ ):

$k$ rows in the set $\mathcal{R}_1^u$ of the form $1^{max} \dots 0$ [rows 1-2 above]
$h + k$ rows (possibly intermingled):
$h$ rows in the set $\mathcal{R}_0^b$ , where $h = \left( \sum_{l=1}^k Z_{l,r} \right) - k$ [rows 3-5 above]
$k$ rows in the set $\mathcal{R}_0^u$ of the form $0 \dots 1$ [rows 6-7 above]

This structure is repeated for each group  $G_i$  with smaller  $N$  value, with the addition of a part consisting of bordered rows arising from splitting blocks of 1's with larger  $N$  values:

$j$ rows in group $G_i$ in the set $\mathcal{R}_1^b$ , where $j = (max - i)k$ [rows 12-13 above for group $G_i = G_2$ ]
--

With reference to Lemma 2.15, we will call a sub-matrix of  $M$  which consists of  $k$  rows in  $\mathcal{R}_1^u$ , having identical  $N$  components, along with their consecutive  $Z$  rotations, an  **$N$ -rep group** – there are two  $N$ -rep groups in Example 4.1: rows 1-7 and 8-11. The repetition of  $N$  blocks gives the  $k$  rows in  $\mathcal{R}_1^u$  at the start of an  $N$ -rep group which yield a run of 0's in the transform  $T$ . Note, however, that rotating  $Z$  suffixes doesn't necessarily give a run of 0's in the transform  $T$  as in Section 2.2, hence we cannot easily determine the number of block pairs here.

We proceed to briefly outline the linear  $B$ -sorting technique for  $B$ -rep words which gives the transform, followed by a description of the inversion method.

The  $B$ -sorting of all conjugates of a  $B$ -rep word is achieved in a similar way to  $B$ -prim words in Section 3, namely extending the definition of  $B$ -order and recording intermediate values during Run Length Encoding. So suppose again that  $\mathbf{u}$  and  $\mathbf{v}$  are suffixes of an input  $\mathbf{x}$  such that  $|\mathbf{u}| < |\mathbf{v}|$ . Then the cases  $\hat{N}_i(\mathbf{u}) \neq \hat{N}_i(\mathbf{v})$ ,  $\hat{N}(\mathbf{u}) = \hat{N}(\mathbf{v}) \neq \varepsilon$ , and empty  $\hat{N}(\mathbf{u})$  or  $\hat{N}(\mathbf{v})$  all follow as before. However, if  $\hat{N}(\mathbf{u})$  is a proper prefix of  $\hat{N}(\mathbf{v})$ , then due to the repetitions of the  $N$  components, we need to order the suffixes using both the  $N$  and  $Z$  values to determine the least – to do this, append the first Lyndon seed of the  $\hat{N}(\mathbf{x})$  repetition together with the adjacent  $\hat{Z}(\mathbf{x})$  values to  $\mathbf{u}$ .

#### 4.1. The Inverse $B$ -BWT - Repetition Case

Rather naturally, the case of  $B$ -rep words with fixed repetition exponent  $k$ , yields the same basic Last First mapping as for  $B$ -prim words but adjusted according to Lemma 2.15.

There are three main steps to the inversion procedure:

1. Compute the array  $IndexRepN_2$  of indexes of the changes in  $N_2$  values – we omit an algorithm for computing this array since it can be achieved



similarly to Algorithm 4 as follows. Scan the transform and count the number of  $N$ -rep groups: they each start with a 0 and end after  $k$  1's. During the scan also count in multiples of  $k$  the other rows, not belonging to  $N$ -rep groups, which are bordered and end with 1. Each  $N$ -rep group starting with  $N > 1$  contributes  $k$  bordered rows to group  $G_1$  - hence we deduce the smallest  $N$  value of the  $N$ -rep groups. Similarly we deduce the second smallest  $N$  value and repeat until we obtain the value of  $N_1 = \max$ . In Example 4.1,  $n = 17$ ,  $k = 2$  and  $T[11]$  is the end of the last  $N$ -rep group. Given 6 bordered rows remaining, and 2  $N$ -rep groups, we see that the first  $N$ -rep group is in group  $G_3$ , and  $\max = 3$ . The rotation of row 6 in  $G_3$ , shows that the first row in  $G_1$ , row 14, must have  $N_2 = 3$ ; and so on.

2. Determine the length of each  $Z$  suffix in the first  $k$  unbordered rows of each  $N$ -rep group – Algorithm 6.
3. Implement the Last First mapping which uses  $IndexRepN_2$  from Step 1, and the array of suffix lengths determined in Step 2 as a look-up table – Algorithm 7.

Algorithm 6 is a linear preprocessing step, bound by constant  $k$  times the length of the longest  $Z$  suffix of each Lyndon seed; Algorithm 7 processes each bit in the transform once with constant work:

**Claim 4.3.** *Procedure  $\text{COMPUTE-LF}_{rep}(T, Z, k, n)$  can be implemented in linear time.*

The reasoning for correctness is similar to that in Claim 3.8, although the  $Z$  suffix of a block pair is obtained from a look-up table, and the current least unselected prefix is from  $k$  bordered or unbordered rows.

**Claim 4.4.** *Procedure  $\text{COMPUTE-LF}_{rep}(T, Z, k, n)$  correctly recovers the input.*

## 5. Bijection

To complete the theory on the new binary transforms, we will establish that they are bijective, that is, if a string  $t = B\text{-BWT}(x)$  for an input  $B$ -word  $x$ , which is specified to be type repetition or primitive, then the appropriate inversion method applied to  $t$  correctly recovers  $x$ .

First we establish bijectivity for the repetition case for a given exponent of repetition  $k$ , and since the primitive case follows similarly, we will see that bijectivity holds for a fixed value of  $k \geq 1$ .

**Lemma 5.1.** *Suppose that  $u, v$  are both  $B$ -rep words with exponent of repetition  $k$ . If  $u \neq v$  then  $B\text{-BWT}(u) \neq B\text{-BWT}(v)$ .*

---

**Algorithm 6** Compute the  $Z$  suffixes for a transform  $T$  based on repetition.

---

```

procedure COMPUTE- $Z_{rep}(T, k, n)$ 
   $Z[1 \dots n]$ ;            $\triangleright$  declare an array initialized to 0
   $i := 1$ ;                  $\triangleright i$  indexes the transform
   $i' := 1$ ;                 $\triangleright i'$  is the start of an  $N$ -rep group

  while  $i < n$  do
     $N := 0$ ;                $\triangleright N$  tracks how many 1's are detected
     $k' := k$ ;                $\triangleright k'$  tracks how many 1's are left to be detected

    while  $N < k$  do        $\triangleright$  find length of the first  $k$   $Z$  suffixes in an  $N$ -rep group
       $\ell := i'$ ;            $\triangleright \ell$  indexes the array  $Z$ 

      for  $j = 1$  to  $k'$  do    $\triangleright$  process suffixes from right to left
        while  $T[\ell]$  is "marked" do  $\ell := \ell + 1$ 
      end while

      if  $T[i] = 0$  then  $Z[\ell] := Z[\ell] + 1$ 
      else "mark"  $T[\ell]$ ;  $N := N + 1$ 
      end if
       $i := i + 1$ ;  $\ell := \ell + 1$ 
    end for
     $k' := k - N$ 
  end while

  while  $T[i] = 1$  do  $i := i + 1$   $\triangleright$  skip bordered rows not in  $N$ -rep groups
  end while

   $i' := i$ 
end while

  return  $Z$ 
end procedure

```

---

---

**Algorithm 7** Compute the Last First Mapping for a transform  $T$  based on repetition.

---

```

procedure COMPUTE-LFrep( $T, Z, k, n$ )
   $currG := 1$ ;           ▷ size of current group
   $prevG := max$ ;         ▷ size of previous group obtained from  $IndexRepN_2$ 
   $p := 10^z$ ;           ▷ initialize the pattern for the matching process;
                        ▷  $z$  is obtained from a linear scan of  $Z[1 \dots k]$ ,
                        ▷ from Algorithm 6, giving the shortest suffix

  while  $|p| < n$  do    ▷ iterate until  $n$  characters have been recovered

     $t := IndexRepN_2[currG, prevG]$ 
    ▷ obtain index of current group in array  $M$  from the array  $IndexRepN_2$ 

    if  $T[t] = 1$  then   ▷ the last bit of the current rotation in  $M$  is 1
       $p := 1 \circ p$ ;      ▷ concatenate the suffix bit 1 to the start of  $p$ 
       $IndexRepN_2[currG, prevG] := IndexRepN_2[currG, prevG] + 1$ ;
       $currG := currG + 1$  ▷ increment the group size

    else ▷ current rotation in  $M$  ends 0 and belongs to an  $N$ -rep group
      if  $t = 1$  then “add marker to start of  $p$ ”   ▷ start of input  $x$ 
      end if
       $p := 10^{Z[t]} \circ p$ ; ▷ concatenate a  $Z$  block of size  $Z[t]$  to the start of  $p$ 
       $IndexRepN_2[currG, prevG] := IndexRepN_2[currG, prevG] + 1$ 
       $prevG := currG$ ;
       $currG := 1$  ▷ the prefix 10 indicates a rotation in group  $G_1$  (Corollary 2.12)
    end if

  end while

  return  $LF_{rep}(x)$  ▷ the input  $x$  is given by the rotation of  $p$  starting with the marker
end procedure

```

---

*Proof.* Suppose that  $B\text{-BWT}(\mathbf{u}) = B\text{-BWT}(\mathbf{v}) = \mathbf{t}$ , then  $|\mathbf{u}| = |\mathbf{v}| = |\mathbf{t}| = n$ . Denote the  $B\text{-BWT}$  matrices of  $\mathbf{u}$ ,  $\mathbf{v}$  as  $M(\mathbf{u})$  and  $M(\mathbf{v})$  respectively. Lemma 2.15 shows that the transform  $\mathbf{t}$  partitions into substrings of two distinct types: those whose elements are the suffixes of  $N$ -rep groups, and those whose elements are the suffixes of bordered rotations formed from splitting  $N$  blocks.

Consider a substring  $\mathbf{s}$  of  $\mathbf{t}$  corresponding to an  $N$ -rep group, which we assume forms rows  $i$  to  $j$  in  $M(\mathbf{u})$  – the first  $N$ -rep group starts at  $\mathbf{t}[1]$  and  $M(\mathbf{u})[1, n]$ . Every such substring  $\mathbf{s}$  starts with  $0^k$  and ends with the  $k^{\text{th}}$  occurrence of a 1 in  $\mathbf{s}$ . So the first  $k$  rows corresponding to this substring  $\mathbf{s}$ ,  $i, i+1, \dots, i+k-1$ , in  $M(\mathbf{u})$  must all end with 0, while the suffix bits of the following  $k$  rows in  $M(\mathbf{u})$  (which must exist), corresponding to  $\mathbf{s}$ , are the  $(n-1)^{\text{th}}$  elements in the rows  $i, i+1, \dots, i+k-1$  in  $M(\mathbf{u})$ . Similarly for the next  $k-d$  rows, where  $d$  is the number of 1's in the previous  $k$  suffix bits. Therefore, there is a 1-1 correspondence between  $\mathbf{s}$  and the  $Z$  block suffixes of rows  $i, i+1, \dots, i+k-1$  in  $M(\mathbf{u})$ . Hence this must also hold for  $M(\mathbf{v})$ . Furthermore, since this is the case for all  $N$ -rep groups, the  $Z$  values in  $\hat{Z}(\mathbf{u})$  and  $\hat{Z}(\mathbf{v})$  must be the same – but not necessarily their order in  $\mathbf{u}$  and  $\mathbf{v}$ .

Applying Step 1 of the inversion method, we first deduce the  $N$  value of the first  $N$ -rep group ( $N \geq 1$ ) and subsequently find the other  $N$  values for the remaining  $N$ -rep groups. Hence we know the  $N$  values of  $\hat{N}(\mathbf{u})$  along with their repetition factor, which similarly must be the same for  $\hat{N}(\mathbf{v})$ , but as above, not necessarily their order. It follows that we also know each of the  $NZ$  pairs (from clockwise rotations of  $N$  suffixes to concatenate with  $Z$  prefixes) which must be the same in both  $\mathbf{u}$  and  $\mathbf{v}$  – but once again, not necessarily their order. So we proceed to show that the order of these block pairs must be the same in both  $\mathbf{u}$  and  $\mathbf{v}$ .

All this implies that both  $\mathbf{u}$  and  $\mathbf{v}$  have the form  $1^{max} \dots 0^{Z_r}$ . This last  $Z$  block suffix is used for starting the Last First mapping, hence the initialized pattern  $p$  is the same in both  $\mathbf{u}$  and  $\mathbf{v}$ .

It also follows that since the  $N$ -rep groups are the same in both matrices, then so to are multiples of  $k$  bordered rows preceding or following these groups. Hence we need to show that the ordering of a bordered group starting with some  $N_i$ , and an  $N$ -rep group also starting with  $N_i$ , is the same in both matrices.

From Step 1 of the inversion method, the entries in the array  $IndexRepN_2$  will be the same for both  $\mathbf{u}$  and  $\mathbf{v}$ . The Last First mapping specifies a sequence of rows visited in the matrices, where the suffix (0 or 1 bit) of a row is the same in both  $M(\mathbf{u})$  and  $M(\mathbf{v})$ . Hence the same sequence of  $N$  values (including split  $N$ 's) is determined for both  $\mathbf{u}$  and  $\mathbf{v}$ . Furthermore, the look-up table given by array  $Z[1 \dots n]$  in Algorithm 6 shows that the  $i^{\text{th}}$   $N$  in both  $\mathbf{u}$  and  $\mathbf{v}$  will have the same adjacent  $Z$  values – that is, the sequence of block pairs is the same in  $\mathbf{u}$  and  $\mathbf{v}$ .  $\square$

Finally, since a similar argument holds for the primitive case, we can state:

**Lemma 5.2.** *Suppose that  $\mathbf{u}$ ,  $\mathbf{v}$  are  $B$ -prim words. If  $\mathbf{u} \neq \mathbf{v}$  then  $B\text{-BWT}(\mathbf{u}) \neq B\text{-BWT}(\mathbf{v})$ .*

Table 1: Improvement gained in the number of blocks for the Rouen Transform versus the Burrows-Wheeler Transform for all the  $B$ -words of length from 4 to 20. For instance there are 30  $B$ -words of length 19 whose Rouen Transform has 8 blocks less than the original Burrows-Wheeler Transform.

	Number of blocks													Total
	8	6	4	2	0	-2	-4	-6	-8	-10	-12	-14	-16	
4				1	2									3
5				2	3	1								6
6				3	6									9
7				7	9	1	1							18
8				13	13	3	1							30
9			2	22	19	12		1						56
10			5	35	42	12	5							99
11			14	61	76	21	12	1	1					186
12		1	29	109	115	62	14	4	1					335
13		4	59	201	219	96	30	19	1	1				630
14		9	133	341	386	182	81	23	6					1161
15		29	260	629	685	347	161	50	20		1			2182
16		74	506	1140	1141	662	300	132	19	5	1			4080
17	2	169	1019	2102	2117	1257	628	236	51	27	1	1		7710
18	7	396	1973	3854	4054	2374	1201	487	148	31	7			14532
19	30	882	3865	7195	7308	4492	2447	940	331	74	28	1	1	27594
20	105	1897	7259	13197	13507	8545	4750	1938	627	226	29	6	1	52377

Interestingly, the  $B$ -prim word 111001101000 (with repetition exponent 1) and the  $B$ -rep word 110101100100 (with exponent 2), both have the same transform, namely 000110010111. This shows the necessity of specifying the type of  $B$ -word, the exponent, and hence the inversion method required: Algorithm 5 or Algorithm 7.

Hence like human twins, these transforms are very similar but distinct: with repetition exponent 1 (although strictly  $> 1$ ), a  $B$ -rep word would require that  $\hat{Z}$  satisfies Lyndon properties, while  $\hat{Z}$  is unrestricted for  $B$ -prim words. Hence  $B$ -prim words can be viewed as a special case of  $B$ -rep words.

## 6. Experimentation

We generated all the  $B$ -words of length from 4 to 20. For each of them we computed both the Rouen (Block Order) Transform and the Burrows-Wheeler Transform; furthermore, we computed the improvement gained in the number of blocks for the Rouen Transform versus the Burrows-Wheeler Transform (see Table 1). These results show that in most cases the two transforms have the same number of blocks; however, in most of the remaining cases the Rouen Transform generates less blocks than the Burrows-Wheeler Transform. This means that a proper decomposition of any given binary word into  $B$ -words could lead to a better compression with the Rouen Transform than with the

Burrows-Wheeler Transform. These encouraging preliminary results show that further experiments need to be done for a thorough classification.

## 7. Conclusion

We have introduced a class of binary Burrows-Wheeler type transforms for binary block order - we believe this is the first approach for transforms designed specifically for binary inputs. We have described two types of  $B$ -words: primitive  $B$ -prim words and repetition  $B$ -rep words; furthermore linear factoring of binary strings into  $B$ -words is established.

Computing both the  $B$ -BWT transforms and inverses are shown to be linear and bijective; the transform applies suffix array techniques while the inversion is based on combinatorial analysis and indexing. In addition to the theoretical contribution, we have demonstrated that it may be worthwhile in practice to implement the Rouen Transform as preprocessing for compression.

An obvious quest for future research is to devise a fully bijective linear transform for binary block order over arbitrary inputs. If the given string is not a  $B$ -word, then it could be factored into these patterned words, however the multi-word approach by Kufleitner [Kuf09] is not readily applicable as it requires sorting the elements of the multi-transform into lexorder in the first column of the matrix  $M$ . Even avoiding specifying the frequency exponent  $k$ , when  $\hat{N}$  is a repetition, with its overhead of  $\log n$ , seems to be a challenge.

Our ultimate interest is to consider degrees of data clustering achievable by different transforms which are derived from a variety of total orders on suitable alphabets. The integer encoding of strings  $\hat{x}$  gives rise to 32 orders by prioritizing  $N$  or  $Z$  values and using  $<$  or  $>$  for the ordering, which includes binary type order [DD03] – as a next step we suggest investigating other new binary transforms to be defined using these orders for the BWT matrices.

In addition to thorough experimental analysis, for instance on genome-based repetitions, we suggest comparing efficiencies obtained when performing a transform on ASCII characters directly, or, converting ASCII characters to binary and then performing a binary transform.

We thus propose that a valuable line of enquiry is to consider binary BWT transforms in the context of bit-optimal compression methods: the goal of the bit-optimal parsing problem is to minimize the length in bits of an encoded text [Lan13].

## 8. Acknowledgements

We thank the anonymous reviewers for carefully reading the manuscript and their helpful comments.

## 9. References

- [ABM08] D. Adjeroh, T. Bell, and A. Mukherjee. *The Burrows-Wheeler Transform: Data Compression, Suffix Arrays, and Pattern Matching*. Springer, 2008.
- [AT05] J. Abel and W. Teahan. Universal text preprocessing for data compression. *IEEE Trans. Comput.*, 54(5):497–507, 2005.
- [BCRS13] M. J. Bauer, A. J. Cox, G. Rosone, and M. Sciortino. Lightweight LCP construction for next-generation sequencing datasets. *CoRR abs/1305.0160*, 2013.
- [BW94] M. Burrows and D. J. Wheeler. A block sorting lossless data compression algorithm. Technical Report 124, Digital Equipment Corporation, 1994.
- [CDP05] M. Crochemore, J. Désarménien, and D. Perrin. A note on the Burrows-Wheeler transformation. *Theoret. Comput. Sci.*, 332(1-3):567–572, 2005.
- [CFL58] K. T. Chen, R. H. Fox, and R. C. Lyndon. Free differential calculus IV — The quotient groups of the lower central series. *Ann. Math.*, 68:81–95, 1958.
- [CGKL13] M. Crochemore, R. Grossi, J. Kärkkäinen, and G. M. Landau. A constant-space comparison-based algorithm for computing the Burrows-Wheeler transform. In *Proc. 24th Annual Symposium on Combinatorial Pattern Matching (CPM)*, pages 74–82, 2013.
- [CHL07] M. Crochemore, C. Hancart, and T. Lecroq. *Algorithms on Strings*. Cambridge University Press, 2007.
- [CT98] B. Chapin and S. R. Tate. Higher compression from the Burrows-Wheeler transform by modified sorting. In *Proc. 1998 Data Compression Conf. (DCC)*, page 532, 1998.
- [Day11] D. E. Daykin. Algorithms for the Lyndon unique maximal factorization. *J. Combin. Math. Combin. Comput.*, 77:65–74, 2011.
- [DD96] T. -N. Danh and D. E. Daykin. The structure of V-order for integer vectors. In A. J. W. Hilton, editor, *Congressus Numerantium*, volume 113, pages 43–53. Utilitas Mat. Pub. Inc., Winnipeg, Canada, 1996.
- [DD03] D. E. Daykin and J. W. Daykin. Lyndon-like and V-order factorizations of strings. *J. Discrete Algorithms*, 1:357–365, 2003.
- [DS14] J. W. Daykin and W. F. Smyth. A bijective variant of the Burrows-Wheeler transform using V-order. *Theoret. Comput. Sci.*, 531:77–89, 2014.

- [Duv83] J.-P. Duval. Factorizing words over an ordered alphabet. *J. Algorithms*, 4(4):363–381, 1983.
- [DW14] J. W. Daykin and B. Watson. A text transformation scheme for degenerate strings. In Costas S. Iliopoulos and Alessio Langiu, editors, *2<sup>nd</sup> International Conference on Algorithms for Big Data (ICABD2014)*, number 1146 in CEUR-WS Proceedings, pages 23–29, Aachen, 2014.
- [FM00] P. Ferragina and G. Manzini. Opportunistic data structures with applications. In *Proc. 41st Annual Symposium on Foundations of Computer Science, (FOCS 2000)*, pages 390–398, 2000.
- [GS12] J. Y. Gil and D. A. Scott. A bijective string sorting transform. *CoRR*, abs/1201.3077, 2012.
- [KA03] P. Ko and S. Aluru. Space efficient linear time construction of suffix arrays. In *Proc. 14th Annual Symposium on Combinatorial Pattern Matching (CPM)*, pages 200–210, 2003.
- [Kär07] J. Kärkkäinen. Fast BWT in small space by blockwise suffix sorting. *Theoret. Comput. Sci.*, 387(3):249–257, 2007.
- [KKP12] J. Kärkkäinen, D. Kempa, and S. J. Puglisi. Slashing the time for BWT inversion. In *Proc. Data Compression Conference (DCC)*, pages 99–108, 2012.
- [KS03] J. Kärkkäinen and P. Sanders. Simple linear work suffix array construction. In *Proc. 30th Internat. Conf. Automata, Languages & Programming*, pages 943–955, 2003.
- [KSB06] J. Kärkkäinen, P. Sanders, and S. Burkhardt. Linear work suffix array construction. *J.ACM*, 53(6):918–936, 2006.
- [Kuf09] M. Kufleitner. On bijective variants of the Burrows-Wheeler transform. In *Proc. Stringology*, pages 65–79, 2009.
- [Lan13] A. Langiu. On parsing optimality for dictionary-based text compression - the Zip case. *J. Discrete Algorithms*, 20:65–70, 2013.
- [Lot83] M. Lothaire. *Combinatorics on Words*. 2nd Edition, Addison-Wesley, Reading, MA (1983); Cambridge University Press, Cambridge (1997). Addison-Wesley, 1983.
- [MRRS05] S. Mantaci, A. Restivo, G. Rosone, and M. Sciortino. An extension of the Burrows Wheeler transform and applications to sequence comparison and data compression. In *Proc. 16th Annual Symposium on Combinatorial Pattern Matching (CPM)*, pages 178–189, 2005.



- [NZC09] G. Nong, S. Zhang, and W. H. Chan. Linear suffix array construction by almost pure induced-sorting. In *Proc. 2009 Data Compression Conf. (DCC)*, pages 193–202, 2009.
- [SLLM09] M. Salson, T. Lecroq, M. Léonard, and L. Mouchard. A four-stage algorithm for updating a Burrows-Wheeler transform. *Theoret. Comput. Sci.*, 410(43):4350–4359, 2009.
- [Smy03] B. Smyth. *Computing Patterns in Strings*. ACM Press Bks. Pearson/Addison-Wesley, 2003.